



University of Jaén Spain

# Sisonddy foururser-eurusik

What is *Textual-Entailment?* 

Given a text t and an hypothesis h we want to define function e which takes these two texts as arguments and returns an answer to the entailment question is h implici *in* t?

YES if h is implicit in t

e(t,h) =

NO otherwise

This can be approximated using Machine Learning algorithms by

$$\hat{e}(t,h) = bc(f,m)$$

where

bc	is a binary classifier
f	is a set of features
т	is the learned model of the
	classifier <i>bc</i>

Therefore, training has to be performed before solving entailment questions which are, with this approach reduced to classification decisions:



# **Combining Lexical-Syntactic Information with Machine** Learning for Recognizing Textual Entailment

A. Montejo-Ráez, J.M. Perea, F. Martínez-Santiago, M.A. García-Cumbreras, M. Martín-Valdivia, A.L. Ureña-López SINAI Research Group – Dpto. de Informática, University of Jaén, 23071 – Jaén (Spain)

<b>Computed feature</b>	5
Lexical similarity	
SIM: Simple Matching	
Semantic distance between each text. When the distance is under a ce are counted as similar.	stem of ertain thre
$SIM_{matching} = \frac{\sum_{i \in H} similarity(i)}{ H }$	where s <sub>L</sub> is the <i>H</i> and <i>T</i>
similarity(i) = $\begin{cases} 1 & \text{if } \exists j \in T, s_L(i, j) > 0.5 \\ 0 & \text{otherwise} \end{cases}$	text sets <i>i,j</i> are tw synsets)
BM: Binary Matching	
Same as previous but:	
similarity(i) = $\begin{cases} 1 & \text{if } \exists j \in T, i = j \\ 0 & \text{otherwise} \end{cases}$	
CSS: Consecutive Subsequence Ma	atching
Counts the number of consecutive appear in both pieces of text.	subseque
Trigrams:	
Same as before, but using words in subsequences of three words.	stead of
Syntactic similarity (SynTree f	eatures
We computed a set of measures from COLLINS from both texts. <i>Aligned</i> terr appearing in both texts), and then, als the parsing, we compute:	the syntans are ide o with the
<ol> <li>Number of aligned terms</li> <li>Number of coincident POS of aligned</li> <li>Number of unmatched POS of aligned</li> <li>Minimal, maximal and average of through the syntactic trees to go from</li> </ol>	ned terms gned term distances om one al
	Computed features Lexical similarity SiM: Simple Matching Semantic distance between each text. When the distance is under a co are counted as similar. $SIM_{mechang} = \frac{\sum_{c \in H} similarity(f)}{ H }$ $similarity(f) = \begin{cases} f \ f \ f \ f \ f \ f \ f \ f \ f \ f$

the hypothesis and the eshold, then both stems

Lin's similarity measure are the hypothesis and of concepts respect. vo concepts (WordNet's

ences of stems that

stems and only on

actic trees obtained with lentified (those stems e POS information of

s differences in nodes ligned term to another

### **Experimental results**

We have performed 6 different experiments, applying combinations of described features and two possible learning algorithms: BBR and TiMBL

SIM	BM	CSS	Trigrams	SynTree	Classifier	Accuracy
Х		X	Х	Х	BBR	0.6475
Х	Х	X	Х	Х	BBR	0.6462
Х		X	Х		BBR	0.6387
Х	Х	X	Х		TiMBL	0.6062
Х		X	Х		TiMBL	0.6037
X		X	X	X	TiMBL	0.5700

# **Conclusions and future work**

Good results areobtained integrating all features, but lexical ones seem to hold better information, being the most crucial feature the simplest one: *trigams* weighting.

We plan to combine the output of several classifiers (like also SVM) to produce a single answer. Also the syntactic information has to be studied in deeper detail. Besides, a model working at semantic level is under implementation.

# References

- Alexander Budanitsky and Graeme Hirst. 200 Semantic distance in WordNet: An experiment application-oriented evaluation of five measur
- Michael Collins. 1999. Head-Driven Statistical Models for Natural Language Parsing. Ph.D. Thesis, University of Pennsylvania
- Walter Daelemans, Jakub Zavrel, Ko van der Sloot and Antal va den Bosc.1998. TiMBL: *Tilburg memory based learner*, version 1.0
- Oscar Ferrandez, Daniel Micolo, Rafael Muño and Manuel Palomar. 2007. Técnicas léxicosintácticas para reconocimiento de implicació textual. Tecnologías de la Información Multilingüe y Multimodal (*in press*)



<ul> <li>Jay J. Jiang and David W. Conrath. 1997. Semantic similarity based on corpus statistics and lexical taxonomy. In Proceedings of International Confeence on Research in Computational Linguistics. Taiwan.</li> <li>Dekang Lin. 1998. An information-theoretic definition of similarity. In Proceedings of the 15<sup>th</sup> International Conference on Machine Learning.</li> <li>Philip Resnik. 1995. Using information content to evaluate semantic similarity. In Proceedings of the 14<sup>th</sup>International Joint Conference on Artificial Intelligence, Montreal.</li> </ul>		
<ul> <li>Dekang Lin. 1998. An information-theoretic definition of similarity. In Proceedings of the 15<sup>th</sup> International Conference on Machine Learning.</li> <li>Philip Resnik. 1995. Using information content to evaluate semantic similarity. In Proceedings of the 14<sup>th</sup>International Joint Conference on Artificial Intelligence, Montreal.</li> </ul>	1. <i>tal,</i> es.	<ul> <li>Jay J. Jiang and David W. Conrath. 1997. Semantic similarity based on corpus statistics and lexical taxonomy. In Proceedings of International Conference on Research in Computational Linguistics. Taiwan.</li> </ul>
<ul> <li>Philip Resnik. 1995. Using information content to evaluate semantic similarity. In Proceedings of the 14<sup>th</sup>International Joint Conference on Artificial Intelligence, Montreal.</li> </ul>		<ul> <li>Dekang Lin. 1998. An information-theoretic definition of similarity. In Proceedings of the 15<sup>th</sup> International Conference on Machine Learning.</li> </ul>
	n	<ul> <li>Philip Resnik. 1995. Using information content to evaluate semantic similarity. In Proceedings of the 14<sup>th</sup>International Joint Conference on Artificial Intelligence, Montreal.</li> </ul>